

Distributed Storage Codes through Hadamard Designs

Dimitris S. Papailiopoulos and Alexandros G. Dimakis

Department of Electrical Engineering

University of Southern California

Los Angeles, CA 90089

Email:{papailio, dimakis}@usc.edu

Abstract—In distributed storage systems that employ erasure coding, the issue of minimizing the total *repair bandwidth* required to exactly regenerate a storage node after a failure arises. This repair bandwidth depends on the structure of the storage code and the repair strategies used to restore the lost data. Minimizing it requires that undesired data during a repair align in the smallest possible spaces, using the concept of interference alignment (IA). Here, a points-on-a-lattice representation of the symbol extension IA of Cadambe *et al.* provides cues to perfect IA instances which we combine with fundamental properties of Hadamard matrices to construct a new storage code with favorable repair properties. Specifically, we build an explicit $(k+2, k)$ storage code over $\mathbb{GF}(3)$, whose single systematic node failures can be repaired with bandwidth that matches exactly the theoretical minimum. Moreover, the repair of single parity node failures generates at most the same repair bandwidth as any systematic node failure. Our code can tolerate any single node failure and any pair of failures that involves at most one systematic failure.

I. INTRODUCTION

The demand for large scale data storage has increased significantly in recent years with applications demanding seamless storage, access, and security for massive amounts of data. When the deployed nodes of a storage network are individually unreliable, as is the case in modern data centers, or peer-to-peer networks, redundancy through erasure coding can be introduced to offer reliability against node failures. However, increased reliability does not come for free: the encoded representation needs to be maintained posterior to node erasures. To maintain the same redundancy when a storage node leaves the system, a new node has to join the array, access some existing nodes, and regenerate the contents of the departed node. This problem is known as the *Code Repair Problem* [3], [1].

The interest in the code repair problem, and specifically in designing repair optimal (n, k) erasure codes, stems from the fact that there exists a fundamental minimum repair bandwidth needed to regenerate a lost node that is substantially less than the size of the encoded data object. MDS erasure storage codes have generated particular interest since they offer maximum reliability for a given storage capacity; such an example is the EvenOdd construction [2]. However, most practical solutions for

storage use existing off-the-shelf erasure codes that are repair inefficient: a single node repair generates network traffic equal to the size of the *entire* stored information.

Designing repair optimal MDS codes, i.e., ones achieving the minimum repair bandwidth bound that was derived in [3], seems to be challenging especially for high rates $\frac{k}{n} \geq \frac{1}{2}$. Recent works by Cadambe *et al.* [11] and Suh *et al.* [12] used the symbol extension IA technique of Cadambe *et al.* [4] to establish the existence, for all n, k , of asymptotically optimal MDS storage codes, that come arbitrarily close to the theoretic minimum repair bandwidth. However, these asymptotic schemes are impractical due to the arbitrarily large file size and field size that they require. Explicit and practical designs for optimal MDS storage codes are constructed roughly for rates $\frac{k}{n} \leq \frac{1}{2}$ [5]–[10], [13], and most of them are based upon the concept of interference alignment. Interestingly, as of now no explicit MDS storage code constructions exist with optimal repair properties for the high data rate regime.¹

Our Contribution: In this work we introduce a new high-rate, explicit, $(k+2, k)$ storage code over $\mathbb{GF}(3)$. Our storage code exploits fundamental properties of Hadamard designs and perfect IA instances pronounced by the use of a lattice representation for the symbol extension IA of Cadambe *et al.* [4]. This representation gives hints for coding structures that allow *exact* instead of asymptotic alignment. Our code exploits these structures and achieves perfect IA without requiring the file size or field size to scale to infinity. Any single systematic node failure can be repaired with bandwidth matching the theoretic minimum and any single parity node failure generates (at most) the same repair bandwidth as any systematic node repair. Our code has two parities but cannot tolerate any two failures: the form presented here can tolerate any single failure and any pair of failures that involves at most one

¹During the submission of this manuscript, two independent works appeared that constructed MDS codes of arbitrary rate that can optimally repair their systematic nodes, see [14], [15].

systematic node	systematic data
1	\mathbf{f}_1
\vdots	\vdots
k	\mathbf{f}_k
parity node	parity data
a	$\mathbf{A}_1^T \mathbf{f}_1 + \dots + \mathbf{A}_k^T \mathbf{f}_k$
b	$\mathbf{B}_1^T \mathbf{f}_1 + \dots + \mathbf{B}_k^T \mathbf{f}_k$

Fig. 1. A $(k+2, k)$ CODED STORAGE ARRAY.

systematic node failure². Here, in contrast to MDS codes, slightly more than k , that is, $k(1 + \frac{1}{2k})$, encoded pieces are required to reconstruct the file object.

II. DISTRIBUTED STORAGE CODES WITH 2 PARITY NODES

In this section, we consider the code repair problem for storage codes with 2 parity nodes. Let a file of size $M = kN$ denoted by the vector $\mathbf{f} \in \mathbb{F}^{kN}$ be partitioned in k parts $\mathbf{f} = [\mathbf{f}_1^T \dots \mathbf{f}_k^T]^T$, each of size N .³ We wish to store this file with rate $\frac{k}{k+2}$ across k systematic and 2 parity storage units each having storage capacity $\frac{M}{k} = N$. To achieve this level of redundancy, the file is encoded using a $(k+2, k)$ distributed storage code. The structure of the storage array is given in Fig. 1, where \mathbf{A}_i and \mathbf{B}_i are $N \times N$ matrices of coding coefficients used by the parity nodes a and b , respectively, to “mix” the contents of the i th file piece \mathbf{f}_i . Observe that the code is in systematic form: k nodes store the k parts of the file and each of the 2 parity nodes stores a linear combination of the k file pieces.

To maintain the same level of redundancy when a node fails or leaves the system, the code repair process has to take place to exactly restore the lost data in a *newcomer* storage component. Let for example a systematic node $i \in \{1, \dots, k\}$ fail. Then, a newcomer joins the storage network, connects to the remaining $k+1$ nodes, and has to download sufficient data to reconstruct \mathbf{f}_i . Observe that the missing piece \mathbf{f}_i exists as a term of a linear combination *only* at each parity node, as seen in Fig. 1. To regenerate it, the newcomer has to download from the parity nodes at least the size of what was lost, i.e., N linearly independent data elements. The downloaded contents from the parity nodes can be represented as a stack of N equations

$$\begin{bmatrix} \mathbf{p}_i^{(a)} \\ \mathbf{p}_i^{(b)} \end{bmatrix} \triangleq \underbrace{\begin{bmatrix} (\mathbf{A}_i \mathbf{V}_i^{(a)})^T \\ (\mathbf{B}_i \mathbf{V}_i^{(b)})^T \end{bmatrix}}_{\text{useful data}} \mathbf{f}_i + \sum_{j=1, j \neq i}^k \underbrace{\begin{bmatrix} (\mathbf{A}_j \mathbf{V}_i^{(a)})^T \\ (\mathbf{B}_j \mathbf{V}_i^{(b)})^T \end{bmatrix}}_{\text{interference by } \mathbf{f}_j} \mathbf{f}_j \quad (1)$$

²Our latest work expands Hadamard designs to construct 2-parity MDS codes that can optimally repair any systematic or parity node failure and m -parity MDS codes that can optimally repair any systematic node failure [16].

³ \mathbb{F} denotes the finite field over which all operations are performed.

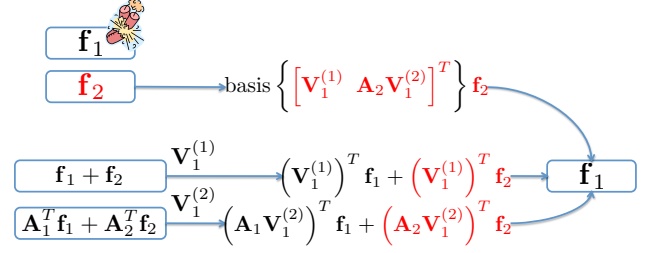


Fig. 2. Repair of a $(4, 2)$ code.

where $\mathbf{p}_i^{(a)}, \mathbf{p}_i^{(b)} \in \mathbb{F}^{\frac{N}{2}}$ are the equations downloaded from parity nodes a and b respectively. Here, $\mathbf{V}_i^{(a)}, \mathbf{V}_i^{(b)} \in \mathbb{F}^{N \times \frac{N}{2}}$ denote the *repair matrices* used to mix the parity contents.⁴ Retrieving \mathbf{f}_i from (II) is equivalent to solving an underdetermined set of N equations in the kN unknowns of \mathbf{f} , with respect to only the N desired unknowns of \mathbf{f}_i . However, this is not possible due to the additive *interference* components that corrupt the desired information in the received equations. These terms are generated by the undesired unknowns $\mathbf{f}_j, j \neq i$, as noted in (II). Additional data need to be downloaded from the systematic nodes, which will “replicate” the interference terms and will be subtracted from the downloaded equations. To erase a single interference term, a download of a basis of equations that generates the corresponding interference term, say $\begin{bmatrix} (\mathbf{A}_s \mathbf{V}_i^{(a)})^T \\ (\mathbf{B}_s \mathbf{V}_i^{(b)})^T \end{bmatrix} \mathbf{f}_j$, suffices. Eventually, when all undesired terms are subtracted, a full rank system of N equations in N unknowns $\begin{bmatrix} (\mathbf{A}_i \mathbf{V}_i^{(a)})^T \\ (\mathbf{B}_i \mathbf{V}_i^{(b)})^T \end{bmatrix} \mathbf{f}_i$ has to be formed. Thus, it can be proven that the *repair bandwidth* to exactly regenerate systematic node i is given by

$$\gamma_i = N + \sum_{j=1, j \neq i}^k \text{rank} \left(\begin{bmatrix} \mathbf{A}_j \mathbf{V}_i^{(a)} & \mathbf{B}_j \mathbf{V}_i^{(b)} \end{bmatrix} \right),$$

where the sum rank term is the aggregate of interference dimensions. Interference alignment plays a key role since the lower the interference dimensions are, the less repair data need to be downloaded. We would like to note that the theoretical minimum repair bandwidth of any node for optimal $(k+2, k)$ MDS codes is exactly $(k+1)\frac{N}{2}$, i.e. half of the remaining contents; this corresponds to each interference spaces having rank $\frac{N}{2}$. This is also true for the systematic parts of non-MDS codes, as long as they have the same problem parameters that were discussed in the beginning of this section, and all the coding matrices have full rank N . An abstract example of a code repair instance for a $(4, 2)$ storage code is given in Fig. 2, where interference terms are marked in red.

To minimize the repair bandwidth γ_i , we need to carefully design both the storage code and the repair matrices.

⁴Here, we consider that the newcomer downloads the same amount of information from both parities. In general this does not need to be the case.

In the following, we provide a 2-parity code that achieves optimal systematic and near optimal parity repair.

III. A NEW STORAGE CODE

We introduce a $(k+2, k)$ storage code over $\mathbb{GF}(3)$, for file sizes $M = k2^k$, with coding matrices

$$\mathbf{A}_i = \mathbf{I}_N, \quad \mathbf{B}_i = \mathbf{X}_i, \quad (2)$$

where $N = 2^k$, $\mathbf{X}_i = \mathbf{I}_{2^{i-1}} \otimes \text{blkdiag}(\mathbf{I}_{\frac{N}{2^i}}, -\mathbf{I}_{\frac{N}{2^i}})$, and $i \in \{1, \dots, k\}$. In Fig. 3, we give the coding matrices of the $(5, 3)$ version of the code.

Theorem 1: The code in (2) has optimally repairable systematic nodes and its parity nodes can be repaired by generating as much repair bandwidth as a systematic repair does. It can tolerate any single node failure, and any pair of failures that contains at most one systematic failure. Moreover, to reconstruct the file at most $k + \frac{1}{2}$ coded blocks are required.

In the following, we present the tools that we use in our derivations. Then, in Sections V and VI we prove Theorem 1.

IV. DOTS-ON-A-LATTICE AND HADAMARD DESIGNS

Optimality during a systematic repair, requires interference spaces collapsing down to the minimum of $\frac{N}{2}$, out of the total N , dimensions. At the same time, useful data equations have to span N dimensions. For the constructions presented here, we consider that the same repair matrix is used by both parities, i.e., $\mathbf{V}_i^{(1)} = \mathbf{V}_i^{(2)} = \mathbf{V}_i$. Hence, for the repair of systematic node $i \in \{1, \dots, k\}$ we optimally require

$$\text{rank}([\mathbf{V}_i \ \mathbf{X}_j \mathbf{V}_i]) = \frac{N}{2}, \quad (3)$$

for all $j \in \{1, \dots, k\} \setminus i$, and at the same time

$$\text{rank}([\mathbf{V}_i \ \mathbf{X}_i \mathbf{V}_i]) = N. \quad (4)$$

The key ingredient of our approach that eventually provides the above is Hadamard matrices.

To motivate our construction, we start by briefly discussing the repair properties of the asymptotic coding schemes of [11], [12]. Consider a 2-parity MDS storage code that requires file sizes $M = k2\Delta^{k-1}$, i.e., $N = 2\Delta^{k-1}$. Its $N \times N$ diagonal coding matrices $\{\mathbf{X}_s\}_{s=1}^k$ have i.i.d. elements drawn uniformly at random from some arbitrarily large finite field \mathbb{F} . During the repair of a systematic node $i \in \{1, \dots, k\}$, the repair matrix \mathbf{V}_i that is used by both parity nodes to mix their contents, has as columns the $\frac{N}{2} = \Delta^{k-1}$ elements of the set

$$\mathcal{V}_i = \left\{ \prod_{s=1, s \neq i}^k \mathbf{X}_s^{x_s} \mathbf{w} : x_s \in \{0, \dots, \Delta-1\} \right\}. \quad (5)$$

Then, we define a map \mathcal{L} from vectors in the set $\left\{ \prod_{s=1}^k \mathbf{X}_s^{x_s} \mathbf{w} : x_s \in \mathbb{Z} \right\}$ to points on the integer lattice

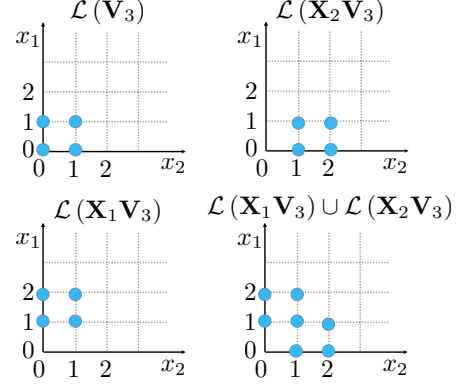


Fig. 4. Here we have $k = 3$, $\frac{N}{2} = 4$, and $\Delta = 2$. Moreover, $\mathcal{L}(\mathbf{V}_3) = \{(0, 0, 0), (0, 1, 0), (1, 0, 0), (1, 1, 0)\}$, $\mathcal{L}(\mathbf{X}_1 \mathbf{V}_3) = \{(1, 0, 0), (1, 1, 0), (2, 0, 0), (2, 1, 0)\}$, and $\mathcal{L}(\mathbf{X}_2 \mathbf{V}_3) = \{(0, 1, 0), (0, 2, 0), (1, 1, 0), (1, 2, 0)\}$.

$\mathbb{Z}^k: \prod_{s=1}^k \mathbf{X}_s^{x_s} \mathbf{w} \xrightarrow{\mathcal{L}} \sum_{s=1}^k x_s \mathbf{e}_s$, where \mathbf{e}_s is the s -th column of \mathbf{I}_{k+1} . Now, consider the induced lattice representation of \mathbf{V}_i

$$\mathcal{L}(\mathbf{V}_i) \triangleq \left\{ \sum_{s=1, s \neq i}^k x_s \mathbf{e}_s; x_s \in \{0, \dots, \Delta-1\} \right\}. \quad (6)$$

Observe that the i -th dimension of the lattice where $\mathcal{L}(\mathbf{V}_i)$ lies on, indicates all possible exponents x_i of \mathbf{X}_i . Then, the products $\mathbf{X}_j \mathbf{V}_i$, $j \neq i$, and $\mathbf{X}_i \mathbf{V}_i$ map to

$$\mathcal{L}(\mathbf{X}_j \mathbf{V}_i) = \left\{ (x_j + 1) \mathbf{e}_j + \sum_{s=1, s \neq j}^k x_s \mathbf{e}_s; x_s \in \{0, \dots, \Delta-1\} \right\}$$

$$\text{and } \mathcal{L}(\mathbf{X}_i \mathbf{V}_i) = \left\{ e_i + \sum_{i=1, s \neq i}^k x_i \mathbf{e}_i; x_s \in \{0, \dots, \Delta-1\} \right\},$$

respectively. In Fig. 2, we give an illustrative example for $k = 3$, and $\Delta = 2$.

Remark 1: Observe how matrix multiplication of \mathbf{X}_i and elements of \mathcal{V}_i manifests itself through the dots-on-a-lattice representation: the product of \mathbf{X}_i with the elements of \mathcal{V}_i shifts the corresponding arrangement of dots along the x_i -axis, i.e., the x_i -coordinate of the initial points gets increased by one.

Asymptotically optimal repair of node i is possible due to the fact that interference spaces asymptotically align

$$\begin{aligned} \frac{\text{rank}([\mathbf{V}_i \ \mathbf{X}_j \mathbf{V}_i])}{\frac{N}{2}} &= \frac{|\mathcal{L}(\mathbf{V}_i) \cup \mathcal{L}(\mathbf{X}_j \mathbf{V}_i)|}{\Delta^{k-1}} \\ &= \frac{|\mathcal{L}(\mathbf{V}_i)| + o(\Delta^{k-1})}{\Delta^{k-1}} \xrightarrow{\Delta \rightarrow \infty} 1, \end{aligned} \quad (7)$$

and useful spaces span N dimensions, that is, $\text{rank}([\mathbf{V}_i \ \mathbf{X}_i \mathbf{V}_i]) = |\mathcal{L}(\mathbf{V}_i) \cup \mathcal{L}(\mathbf{X}_i \mathbf{V}_i)| = 2\Delta^{k-1}$, with arbitrarily high probability for sufficiently large field sizes.

The question that we answer here is the following: How can we design the coding and the repair matrices such that *i)* exact interference alignment is possible and *ii)* the full

$$\mathbf{X}_1 = \text{diag} \left(\begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix} \right), \quad \mathbf{X}_2 = \text{diag} \left(\begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \\ 1 \\ -1 \end{bmatrix} \right), \quad \mathbf{X}_3 = \text{diag} \left(\begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \\ 1 \\ -1 \end{bmatrix} \right)$$

Fig. 3. The coding matrices of a repair optimal $(5, 3)$ code over $\mathbb{GF}(3)$.

rank property is satisfied, for fixed in k file size and field size? We first address the first part. We want to design the code such that the space of the repair matrix is invariant to any transformation by matrices generating its columns, i.e., $\mathcal{L}(\mathbf{X}_j \mathbf{V}_i) = \mathcal{L}(\mathbf{V}_i)$. This is possible when

$$\begin{aligned} \mathcal{L}(\mathbf{X}_j \mathbf{V}_i) &= \left\{ (x_j + 1)\mathbf{e}_j + \sum_{s=1, s \neq j}^k x_s \mathbf{e}_s; x_s \in \{0, \dots, \Delta - 1\} \right\} \\ &= \left\{ x_j \mathbf{e}_j + \sum_{s=1, s \neq j}^k x_s \mathbf{e}_s; x_s \in \{0, \dots, \Delta - 1\} \right\} = \mathcal{L}(\mathbf{V}_i), \end{aligned}$$

that is, when the matrix powers “wrap around” upon reaching their modulus Δ . This wrap-around property is obtained when the diagonal coding matrices have elements that are roots of unity.

Lemma 1: For diagonal matrices, $\mathbf{X}_1, \dots, \mathbf{X}_k$, whose elements are Δ -th roots of unity, i.e., $\mathbf{X}_s^\Delta = \mathbf{X}_s^0$, for all $s \in \{1, \dots, k\}$, we have that $\mathcal{L}(\mathbf{X}_j \mathbf{V}_i) = \mathcal{L}(\mathbf{V}_i)$, for all $i \in \{1, \dots, k\} \setminus j$.

However, arbitrary diagonal matrices whose elements are roots of unity are not sufficient to ensure the full rank property of the useful data repair space $[\mathbf{V}_i \ \mathbf{X}_i \mathbf{V}_i]$. In the following we prove that the full rank property along with perfect IA is guaranteed when we set $N = 2^k$, $\mathbf{X}_i = \mathbf{I}_{2^{i-1}} \otimes \text{blkdiag}(\mathbf{I}_{\frac{N}{2^i}}, -\mathbf{I}_{\frac{N}{2^i}})$, and consider the set

$$\mathcal{H}_N = \left\{ \prod_{i=1}^k \mathbf{X}_i^{x_i} \mathbf{w} : x_i \in \{0, 1\} \right\}. \quad (8)$$

Interestingly, there is a one-to-one correspondence between the elements of \mathcal{H}_N and the columns of a Hadamard matrix.

Lemma 2: Let an $N \times N$ Hadamard matrix of the Sylvester’s construction

$$\mathbf{H}_N \triangleq \begin{bmatrix} \mathbf{H}_{\frac{N}{2}} & \mathbf{H}_{\frac{N}{2}} \\ \mathbf{H}_{\frac{N}{2}} & -\mathbf{H}_{\frac{N}{2}} \end{bmatrix}, \quad (9)$$

with $\mathbf{H}_1 = 1$. Then, \mathbf{H}_N is full-rank with mutually orthogonal columns, that are the N elements of \mathcal{H}_N . Moreover, any two columns of \mathbf{H}_N differ in $\frac{N}{2}$ positions.

The proof is omitted due to lack of space. To illustrate the connection between \mathcal{H}_N and \mathbf{H}_N we “decompose” the Hadamard matrix of order 4

$$\mathbf{H}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix} = [\mathbf{w} \ \mathbf{X}_2 \mathbf{w} \ \mathbf{X}_1 \mathbf{w} \ \mathbf{X}_2 \mathbf{X}_1 \mathbf{w}], \quad (10)$$

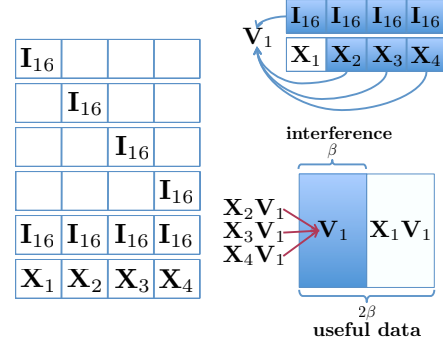


Fig. 5. The coding matrices of our $(6, 4)$ code are given. We illustrate the “absorbing” properties of the repair matrix for systematic node 1. The column space of the repair matrices is invariant to the corresponding blue blocks. This results in interference spaces aligning in exactly half of the dimensions available.

where $\mathbf{X}_1 = \text{diag} \left(\begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix} \right)$ and $\mathbf{X}_2 = \text{diag} \left(\begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \end{bmatrix} \right)$. Due to the commutativity of \mathbf{X}_1 and \mathbf{X}_2 , the columns of \mathbf{H}_4 are also the elements of $\mathcal{H}_4 = \{\mathbf{w}, \mathbf{X}_1 \mathbf{w}, \mathbf{X}_2 \mathbf{w}, \mathbf{X}_1 \mathbf{X}_2 \mathbf{w}\}$.

By using \mathcal{H}_N as our “base” set, we are able to obtain perfect alignment condition due to the wrap around property of its elements; the full rank condition will be also satisfied due to the mutual orthogonality of these elements.

V. REPAIRING SINGLE NODE FAILURES

A. Systematic Repairs

Let systematic node $i \in \{1, \dots, k\}$ fail. Then, we pick the columns of the repair matrix as a set of $\frac{N}{2}$ vectors whose lattice representation is invariant to all \mathbf{X}_j s but to one key matrix \mathbf{X}_i . We specifically construct the $N \times \frac{N}{2}$ repair matrix \mathbf{V}_i whose columns have a one-to-one correspondence with the elements of the set

$$\mathcal{V}_i = \left\{ \prod_{s=1, s \neq i}^k \mathbf{X}_s^{x_s} \mathbf{w} : x_s \in \{0, 1\} \right\}. \quad (11)$$

First, observe that \mathbf{V}_i is full column rank since it is a collection of $\frac{N}{2}$ distinct columns from \mathcal{H}_N . Then, we have the following lemma.

Lemma 3: For any $i, j \in \{1, 2, \dots, k\}$, we have that

$$\begin{aligned} \text{rank}([\mathbf{V}_i \ \mathbf{X}_j \mathbf{V}_i]) &= |\mathcal{L}(\mathbf{V}_i) \cup \mathcal{L}(\mathbf{X}_j \mathbf{V}_i)| \\ &= \begin{cases} N, & i = j \\ \frac{N}{2}, & i \neq j \end{cases}. \end{aligned} \quad (12)$$

The above holds due to each element of \mathcal{H}_N being associated with a unique power tuple. Then, the columns of $[\mathbf{V}_i \mathbf{X}_i \mathbf{V}_i]$ are exactly the elements of \mathcal{H}_N , since

$$\begin{aligned} \mathcal{L}(\mathbf{V}_i) \cup \mathcal{L}(\mathbf{X}_i \mathbf{V}_i) &= \left\{ \sum_{s=1, s \neq i}^k x_i \mathbf{e}_i; x_i \in \{0, 1\} \right\} \\ &\cup \left\{ \mathbf{e}_i + \sum_{s=1, s \neq i}^k x_i \mathbf{e}_i; x_i \in \{0, 1\} \right\} \\ &= \mathcal{L}(\mathbf{H}_N). \end{aligned} \quad (13)$$

Moreover, the set of columns in \mathbf{V}_i are identical to the set of columns of $\mathbf{X}_j \mathbf{V}_i$, i.e., $\mathcal{L}(\mathbf{V}_i) = \mathcal{L}(\mathbf{X}_j \mathbf{V}_i)$, for $j \neq i$, due to Lemmata 1 and 2. Therefore, the interference spaces span $\frac{N}{2}$ dimensions, which is the theoretic minimum, and the desired data space during any systematic node repair is full-rank, since it has as columns all columns of \mathbf{H}_N .

Hence, we conclude that a single systematic node of the code can be repaired with bandwidth $(k+1)\frac{N}{2} = \frac{k+1}{2k}M$. In Fig. 4, we depict a $(6, 4)$ code of our construction, along with the illustration of the repair spaces.

B. Parity repairs

Here, we prove that a single parity node repair generates at most the repair bandwidth of a single systematic repair. Let parity node a fail. Then, observe that if the newcomer uses the $N \times N$ repair matrix $\mathbf{V}_a^{(b)} = \mathbf{X}_1$ to multiply the contents of parity node b , then it downloads $\mathbf{X}_1 \left(\sum_{i=1}^k \mathbf{X}_i \mathbf{f}_i \right) = \mathbf{f}_1 + \sum_{i=2}^k \mathbf{X}_1 \mathbf{X}_i \mathbf{f}_i$. Observe, that the component corresponding to systematic part \mathbf{f}_1 appears the same in the linear combination stored at the lost parity. By Lemma 2, each of the remaining blocks, $\mathbf{X}_1 \mathbf{X}_i \mathbf{f}_i$ share exactly $\frac{N}{2}$ indices with equal elements to the same $\frac{N}{2}$ indices of $\mathbf{X}_i \mathbf{f}_i$ which was lost, for any $i \in \{2, \dots, k\}$. This is due to the fact that the diagonal elements of matrices $\mathbf{X}_1 \mathbf{X}_i$ and \mathbf{X}_i are the elements of some two columns of \mathbf{H}_N . Therefore, the newcomer has to download from systematic node $j \in \{2, \dots, k\}$, the $\frac{N}{2}$ entries that parity a 's component $\mathbf{X}_j \mathbf{f}_j$ differs from the term $\mathbf{X}_1 \mathbf{X}_j \mathbf{f}_j$ of the downloaded linear combination. Hence, the first parity can be repaired with bandwidth at most $N + (k-1)\frac{N}{2} = (k+1)\frac{N}{2}$.⁵ The repair of parity node b can be performed in the same manner.

VI. ERASURE RESILIENCY

Our code can tolerate any single node failure and any two failures with at most one of them being a systematic one. A double systematic and parity node failure can be treated by first reconstructing the lost systematic node from the remaining parity, and then reconstructing the lost parity from all the systematic nodes. However, two simultaneous systematic node failures cannot be tolerated. Consider for example the corresponding matrix when we

connect to nodes $\{1, \dots, k-2\}$ and both parities:

$$\begin{bmatrix} \mathbf{I}_N & \dots & \mathbf{0}_{N \times N} & \mathbf{0}_{N \times N} & \mathbf{0}_{N \times N} \\ \vdots & & \vdots & \vdots & \vdots \\ \mathbf{0}_{N \times N} & \dots & \mathbf{I}_N & \mathbf{0}_{N \times N} & \mathbf{0}_{N \times N} \\ \mathbf{I}_N & \dots & \mathbf{I}_N & \mathbf{I}_N & \mathbf{I}_N \\ \mathbf{X}_1 & \dots & \mathbf{X}_{k-2} & \mathbf{X}_{k-1} & \mathbf{X}_k \end{bmatrix} \mathbf{f}. \quad (14)$$

The rank of this $kN \times kN$ matrix is $(k-1)N + \frac{N}{2}$ due to the submatrix $\begin{bmatrix} \mathbf{I}_N & \mathbf{I}_N \\ \mathbf{X}_{k-1} & \mathbf{X}_k \end{bmatrix}$ having rank $\frac{3N}{2}$. For these cases, an extra download of $\frac{N}{2}$ equations is required to decode the file, i.e., an aggregate download of $kN + \frac{N}{2}$ equations, or $k + \frac{1}{2}$ encoded pieces.

REFERENCES

- [1] The Coding for Distributed Storage wiki <http://tinyurl.com/storagecoding>
- [2] M. Blaum, J. Brady, J. Bruck, and J. Menon, "EVENODD: An efficient scheme for tolerating double disk failures in raid architectures," in *IEEE Trans. on Computers*, 1995.
- [3] A. G. Dimakis, P. G. Godfrey, Y. Wu, M. J. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems," in *IEEE Trans. on Inform. Theory*, vol. 56, pp. 4539 – 4551, Sep. 2010.
- [4] V. R. Cadambe and S. A. Jafar, "Interference alignment and the degrees of freedom for the K user interference channel," *IEEE Trans. on Inform. Theory*, vol. 54, pp. 3425–3441, Aug. 2008.
- [5] Y. Wu and A. G. Dimakis, "Reducing repair traffic for erasure coding-based storage via interference alignment," in *Proc. IEEE Int. Symp. on Information Theory (ISIT)*, Seoul, Korea, Jul. 2009.
- [6] D. Cullina, A. G. Dimakis, and T. Ho, "Searching for minimum storage regenerating codes," in *Allerton Conf. on Control, Comp., and Comm.*, Urbana-Champaign, IL, September 2009.
- [7] K.V. Rashmi, N. B. Shah, P. V. Kumar, and K. Ramchandran "Exact regenerating codes for distributed storage," in *Allerton Conf. on Control, Comp., and Comm.*, Urbana-Champaign, IL, September 2009.
- [8] N. B. Shah, K. V. Rashmi, P. V. Kumar, and K. Ramchandran, "Explicit codes minimizing repair bandwidth for distributed storage," in *Proc. IEEE ITW*, Jan. 2010.
- [9] C. Suh and K. Ramchandran, "Exact regeneration codes for distributed storage repair using interference alignment," in *Proc. 2010 IEEE Int. Symp. on Inform. Theory (ISIT)*, Seoul, Korea, Jun. 2010.
- [10] Y. Wu. "A construction of systematic MDS codes with minimum repair bandwidth." Submitted to *IEEE Transactions on Information Theory*, Aug. 2009. Preprint available at <http://arxiv.org/abs/0910.2486>.
- [11] V. Cadambe, S. Jafar, and H. Maleki, "Distributed data storage with minimum storage regenerating codes - exact and functional repair are asymptotically equally efficient," in *2010 IEEE Intern. Workshop on Wireless Network Coding (WiNC)*, Apr. 2010.
- [12] C. Suh and K. Ramchandran, "On the existence of optimal exact-repair MDS codes for distributed storage," Apr. 2010. Preprint available online at <http://arxiv.org/abs/1004.4663>
- [13] K. Rashmi, N. B. Shah, and P. V. Kumar, "Optimal exact-regenerating codes for distributed storage at the MSR and MBR points via a product-matrix construction," submitted to *IEEE Transactions on Information Theory*, Preprint available online at <http://arxiv.org/pdf/1005.4178>.
- [14] I. Tamo, Z. Wang, and J. Bruck "MDS Array Codes with Optimal Rebuilding," to appear at *ISIT 2011*, preprint available at <http://arxiv.org/abs/1103.3737>
- [15] V. R. Cadambe, C. Huang, and J. Li, "Permutation codes: optimal exact-repair of a single failed node in MDS code based distributed storage systems," to appear at *ISIT 2011*, preprint available at <http://newport.eecs.uci.edu/~vcadambe/permutations.pdf>
- [16] D. S. Papailiopoulos, A. G. Dimakis, and V. R. Cadambe, "Repair optimal erasure codes through hadamard designs," preprint available at <http://www-scf.usc.edu/~papailio/>

⁵By "at most" we mean that this result is proved using an achievable scheme, however, we do not prove that it is optimal.